

Claims:

What is claimed is:

1. A computerized method of workload balancing for improved availability within a multitude of applications-servers and a multitude of application-clients interconnected with said application-servers by a communication network, said method comprising
 - 6 caching within an application-client, availability data of a subset of currently active
7 application-servers as potential target application-servers, and
8 selecting by execution of an application-request by said application-client, an
9 application-server from said subset based on a load-balancing decision of said
10 application-client as target application-server and is sending said application-request to said
11 target application-server.
 2. A method of workload balancing according to claim 1 further comprising,
 - 12 retrieving by said application-client, said subset of currently active application-servers from
13 a cluster database, said cluster database storing availability data of currently active
14 application-servers.
 3. A method of workload balancing according to claim 2 further comprising,
 - 15 choosing said subset of currently active application-servers by approximating an even
16 distribution of said application-servers with respect to other subsets of other
17 application-clients.
 4. A method of workload balancing according to claim 2 further comprising,

retrieving said subset of currently active application-servers by issuing a retrieve request to a server-monitor, said server-monitor monitoring said application-servers activity status.

5. A method of workload balancing according to claim 4 further comprising,

choosing by said server-monitor, said subset of currently active application-servers from said cluster database by a random selection.

6. A method of workload balancing according to claim 4 further comprising,

choosing by said server-monitor, said subset of currently active application-servers from said cluster database based on a load-balancing decision of said server-monitor.

7. A method of workload balancing according to claim 6 further comprising,

choosing by said server-monitor, said subset of currently active application-servers by approximating an even distribution of said application-servers with respect to other subsets of other application-clients.

8. A method of workload balancing according to claim 4 further comprising,

sending by said application-client together with said retrieve request, an indication of the expected amount of workload to be generated by said application-client, and

basing said server-monitor's load-balancing decision on said indication of the expected amount of workload and/or said application-servers processing power and/or said application-servers work load.

9. A method of workload balancing according to claim 1,

2 wherein selecting by said application-client comprises selecting said target
3 application-server randomly.

1 10. A method of workload balancing according to claim 1,

2 wherein said availability data of said subset of currently active application-servers
3 comprises said application-servers processing power and in addition said
4 application-servers work load and/or said application-servers observed response time and
5 wherein selecting by said application-client comprises selecting said target
6 application-server approximating an even distribution of workload within said subset of
7 currently active application-servers.

11. A method of workload balancing according to claim 1 further comprising,

dynamically updating said availability data of said subset of currently active
application-servers by explicitly requesting from said server-monitor an update of said
subset of currently active application-servers, and/or

dynamically updating by said server-monitor, said availability data of said subset of
currently active application-servers by receiving an update of said subset of currently active
application-servers.

12. A method of workload balancing according to claim 11,

2 wherein said explicitly requesting from said server-monitor an update of said subset of
3 currently active application-servers occurs at one of:

4 prior to an application-request, or
5 after sending an application-request, or
6 periodically, or
7 randomly.

1 13. A method of workload balancing according to claim 11,

2 wherein, in the event an additional application-server has become active or if an active
3 application-server has become inactive, said server-monitor takes the initiative of sending
4 said update of said subset of currently active application-servers.

1 14. A method of workload balancing according to claim 13,

2 wherein said server-monitor sending said update of said subset of currently active
3 application-servers comprises a lazy-multi-cast including;

4 iteratively selecting with a time delay, a collection of application-clients for said multitude
5 of application-clients, and

6 sending each application-client of said collection an update of said subse of currently active
7 application-servers.

8 15. A method of workload balancing according to claim 4,

1 wherein said application-servers each comprise a hot pool of one or a multitude of
2 application-instances, said application-servers being executed in the same or a multitude of
3 address-spaces, and

4 wherein for each of said hot pools a watchdog is monitoring said hot pool's activity status,
5 and

6 wherein each of said watchdogs is informing said server-monitor on the workload of said
7 monitored application-servers and its current activity status.
8

- 1 16. A method of workload balancing according to claim 15,
2 wherein at least one of said watchdogs is monitoring the activity status of other watchdogs
3 and, in the event that said one watchdog detects the unavailability of any one of said other
4 watchdogs, said one watchdog informs said server-monitor of said corresponding activity
5 status and said server-monitor updates said cluster database accordingly.
- 1 17. An apparatus for workload balancing giving improved availability within a multitude of
2 applications-servers and a multitude of application-clients interconnected with said
3 application-servers by a communication network, said apparatus comprising
4 means for caching within an application-client, availability data of a subset of currently
5 active application-servers as potential target application-servers, and
6
7 means for selecting by execution of an application-request by said application-client, an
8 application-server from said subset based on a load-balancing decision of said
9 application-client as target application-server and is sending said application-request to said
10 target application-server.
- 1 18. An apparatus for workload balancing according to claim 17 further comprising,
2 means for retrieving by said application-client, said subset of currently active
3 application-servers from a cluster database, said cluster database storing availability data of
4 currently active application-servers.
- 1 19. An apparatus for workload balancing according to claim 18 further comprising,
2 means for choosing said subset of currently active application-servers by approximating an
3 even distribution of said application-servers with respect to other subsets of other
4 application-clients.

- 1 20. An apparatus for workload balancing according to claim 18 further comprising,
2 means for retrieving said subset of currently active application-servers by issuing a retrieve
3 request to a server-monitor, said server-monitor monitoring said application-servers
4 activity status.
- 1 21. An apparatus for workload balancing according to claim 20 further comprising,
2 means for choosing by said server-monitor, said subset of currently active
3 application-servers from said cluster database by a random selection.
22. An apparatus for workload balancing according to claim 20,
means for choosing by said server-monitor, said subset of currently active
application-servers from said cluster database based on a load-balancing decision of said
server-monitor.
23. An apparatus for workload balancing according to claim 22 further comprising,
means for choosing by said server-monitor, said subset of currently active
application-servers by approximating an even distribution of said application-servers with
respect to other subsets of other application-clients.
- 1 24. An apparatus for workload balancing according to claim 20 further comprising,
2 means for sending by said application-client together with said retrieve request, an
3 indication of the expected amount of workload to be generated by said application-client,
4 and

means for basing said server-monitor's load-balancing decision on said indication of the expected amount of workload and/or said application-servers processing power and/or said application-servers work load.

25. An apparatus for workload balancing according to claim 17,

wherein said means for selecting by said application-client comprises means for selecting said target application-server randomly.

26. An apparatus for workload balancing according to claim 17,

wherein said availability data of said subset of currently active application-servers comprises said application-servers processing power and in addition said application-servers work load and/or said application-servers observed response time and wherein means for selecting by said application-client comprises means for selecting said target application-server approximating an even distribution of workload within said subset of currently active application-servers.

27. An apparatus for workload balancing according to claim 17 further comprising,

means for dynamically updating said availability data of said subset of currently active application-servers by explicitly requesting from said server-monitor an update of said subset of currently active application-servers, and/or

means for dynamically updating by said server-monitor, said availability data of said subset of currently active application-servers by receiving an update of said subset of currently active application-servers.

28. An apparatus for workload balancing according to claim 27,

2 wherein said explicitly requesting from said server-monitor an update of said subset of
3 currently active application-servers occurs at one of:

4 prior to an application-request, or
5 after sending an application-request, or
6 periodically, or
7 randomly.

1 29. An apparatus for workload balancing according to claim 27,

2 wherein, in the event an additional application-server has become active or if an active
3 application-server has become inactive, said server-monitor takes the initiative of sending
4 said update of said subset of currently active application-servers.

5 30. An apparatus for workload balancing according to claim 29,

6 wherein said server-monitor sending said update of said subset of currently active
7 application-servers comprises a lazy-multi-cast including;

means for iteratively selecting with a time delay, a collection of application-clients for said
multitude of application-clients, and

means for sending each application-client of said collection an update of said subse of
currently active application-servers.

1 31. An apparatus for workload balancing according to claim 20,

2 wherein said application-servers each comprise a hot pool of one or a multitude of
3 application-instances, said application-servers being executed in the same or a multitude of
4 address-spaces, and

5 wherein for each of said hot pools a watchdog is monitoring said hot pool's activity status,
6 and

7 wherein each of said watchdogs is informing said server-monitor on the workload of said
8 monitored application-servers and its current activity status.

1 32. An apparatus for workload balancing according to claim 31,

2 wherein at least one of said watchdogs is monitoring the activity status of other watchdogs
3 and, in the event that said one watchdog detects the unavailability of any one of said other
4 watchdogs, said one watchdog informs said server-monitor of said corresponding activity
5 status and said server-monitor updates said cluster database accordingly.

6 33. A computer program product comprising a computer usable medium having computer
7 readable program code means therein for workload balancing giving improved availability
8 within a multitude of applications-servers and a multitude of application-clients
9 interconnected with said application-servers by a communication network, said computer
10 program product comprising

11 computer readable program code means for caching within an application-client,
12 availability data of a subset of currently active application-servers as potential target
1 application-servers, and

2 computer readable program code means for selecting by execution of an
3 application-request by said application-client, an application-server from said subset based
4 on a load-balancing decision of said application-client as target application-server and is
5 sending said application-request to said target application-server.

6 34. The computer program product for workload balancing according to claim 33 further
7 comprising,

computer readable program code means for retrieving by said application-client, said subset of currently active application-servers from a cluster database, said cluster database storing availability data of currently active application-servers.

35. The computer program product for workload balancing according to claim 34 further comprising,

computer readable program code means for choosing said subset of currently active application-servers by approximating an even distribution of said application-servers with respect to other subsets of other application-clients.

36. The computer program product for workload balancing according to claim 34 further comprising,

computer readable program code means for retrieving said subset of currently active application-servers by issuing a retrieve request to a server-monitor, said server-monitor monitoring said application-servers activity status.

37. The computer program product for workload balancing according to claim 36 further comprising,

computer readable program code means for choosing by said server-monitor, said subset of currently active application-servers from said cluster database by a random selection.

38. The computer program product for workload balancing according to claim 36,

computer readable program code means for choosing by said server-monitor, said subset of currently active application-servers from said cluster database based on a load-balancing decision of said server-monitor.

39. The computer program product for workload balancing according to claim 38 further comprising,

computer readable program code means for choosing by said server-monitor, said subset of currently active application-servers by approximating an even distribution of said application-servers with respect to other subsets of other application-clients.

40. The computer program product for workload balancing according to claim 36 further comprising,

computer readable program code means for sending by said application-client together with said retrieve request, an indication of the expected amount of workload to be generated by said application-client, and

computer readable program code means for basing said server-monitor's load-balancing decision on said indication of the expected amount of workload and/or said application-servers processing power and/or said application-servers work load.

41. The computer program product for workload balancing according to claim 33,

wherein said computer readable program code means for selecting by said application-client comprises computer readable program code means for selecting said target application-server randomly.

42. The computer program product for workload balancing according to claim 33,

wherein said availability data of said subset of currently active application-servers comprises said application-servers processing power and in addition said application-servers work load and/or said application-servers observed response time and

wherein computer readable program code means for selecting by said application-client comprises computer readable program code means for selecting said target application-server approximating an even distribution of workload within said subset of currently active application-servers.

43. The computer program product for workload balancing according to claim 33 further comprising,

computer readable program code means for dynamically updating said availability data of said subset of currently active application-servers by explicitly requesting from said server-monitor an update of said subset of currently active application-servers, and/or

computer readable program code means for dynamically updating by said server-monitor, said availability data of said subset of currently active application-servers by receiving an update of said subset of currently active application-servers.

44. The computer program product for workload balancing according to claim 43,

wherein said explicitly requesting from said server-monitor an update of said subset of currently active application-servers occurs at one of:

prior to an application-request, or
after sending an application-request, or
periodically, or
randomly.

45. The computer program product for workload balancing according to claim 43,

wherein, in the event an additional application-server has become active or if an active application-server has become inactive, said server-monitor takes the initiative of sending said update of said subset of currently active application-servers.

- 1
- 2
- 3
- 4
- 5
- 6
- 7

1
2
3
4
5
6
7
8

- 1
- 2
- 3
- 4
- 5